



Dynamic Memory Management

Professor Jennifer Rexford
<http://www.cs.princeton.edu/~jrex>

1



Goals of Today's Lecture

- **Dynamic memory management**
 - Garbage collection by the run-time system (Java)
 - Manual deallocation by the programmer (C, C++)
- **Challenges of manual deallocation**
 - Arbitrary request sizes in an arbitrary order
 - Complex evolution of heap as a program runs
- **Design decisions for the “K&R” implementation**
 - Circular linked-list of free blocks with a “first fit” allocation
 - Coalescing of adjacent blocks to create larger blocks

2

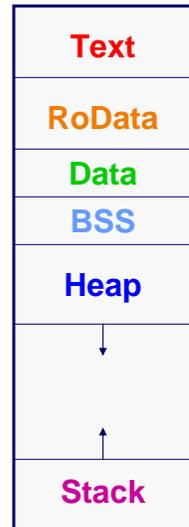
Memory Layout: Heap



```
char* string = "hello";
int iSize;

char* f()
{
    char* p;
    scanf("%d", &iSize);
    p = malloc(iSize);
    return p;
}
```

Needed when required memory size is not known before the program runs



3

Allocating & Deallocating Memory



- **Dynamically allocating memory**
 - Programmer explicitly requests space in memory
 - Space is allocated dynamically on the heap
 - E.g., using “malloc” in C, and “new” in Java
- **Dynamically deallocating memory**
 - Must reclaim or recycle memory that is never used again
 - To avoid (eventually) running out of memory
- **“Garbage”**
 - Allocated block in heap that will not be accessed again
 - Can be reclaimed for later use by the program

4

Option #1: Garbage Collection



- *Run-time system does garbage collection (Java)*
 - Automatically determines objects that can't be accessed
 - And then reclaims the resources used by these objects

```
Object x = new Foo();
Object y = new Bar();
x = new Quux();

if (x.check_something()) {
    x.do_something(y);
}
System.exit(0);
```

Object Foo() is never used again!

5

Challenges of Garbage Collection



- *Detecting the garbage is not always easy*
 - “if (complex_function(y)) x = Quux();”
 - Run-time system cannot collect *all* of the garbage
- *Detecting the garbage introduces overhead*
 - Keeping track of references to objects (e.g., counter)
 - Scanning through accessible objects to identify garbage
 - Sometimes walking through a large amount of memory
- *Cleaning the garbage leads to bursty delays*
 - E.g., periodic scans of the objects to hunt for garbage
 - Leading to unpredictable “freeze” of the running program
 - Very problematic for real-time applications
 - ... though good run-time systems avoid long freezes

6

Option #2: Manual Deallocation



- *Programmer deallocates the memory (C and C++)*
 - Manually determines which objects can't be accessed
 - And then explicitly returns the resources to the heap
 - E.g., using "free" in C or "delete" in C++
- **Advantages**
 - Lower overhead
 - No unexpected "pauses"
 - More efficient use of memory
- **Disadvantages**
 - More complex for the programmer
 - Subtle memory-related bugs
 - Security vulnerabilities in the (buggy) code

7

Manual Deallocation Can Lead to Bugs



- **Dangling pointers**
 - Programmer frees a region of memory
 - ... but still has a pointer to it
 - Dereferencing pointer reads or writes *nonsense values*

```
int main(void) {  
    char *p;  
    p = malloc(10);  
    ...  
    free(p);  
    ...  
    putchar(*p);  
}
```

May print
nonsense
character.

8

Manual Deallocation Can Lead to Bugs

- Memory leak

- Programmer neglects to free unused region of memory
- So, the space can never be allocated again
- Eventually may consume all of the available memory

```
void f(void) {  
    char *s;  
    s = malloc(50);  
    return;  
}  
  
int main(void) {  
    while (1) f();  
    return 0;  
}
```

Eventually,
malloc()
returns NULL

9

Manual Deallocation Can Lead to Bugs

- Double free

- Programmer mistakenly frees a region more than once
- Leading to corruption of the heap data structure
- ... or premature destruction of a *different* object

```
int main(void) {  
    char *p, *q;  
    p = malloc(10);  
    ...  
    free(p);  
    q = malloc(10);  
    free(p);  
    ...  
}
```

Might free the
space allocated
to **q**!

10

Challenges for Malloc and Free



- Malloc() may ask for an arbitrary number of bytes
- Memory may be allocated & freed in different order
- Cannot reorder requests to improve performance

```
char *p1 = malloc(3);
char *p2 = malloc(1);
char *p3 = malloc(4);
free(p2);
char *p4 = malloc(6);
free(p3);
char *p5 = malloc(2);
free(p1);
free(p4);
free(p5);
```

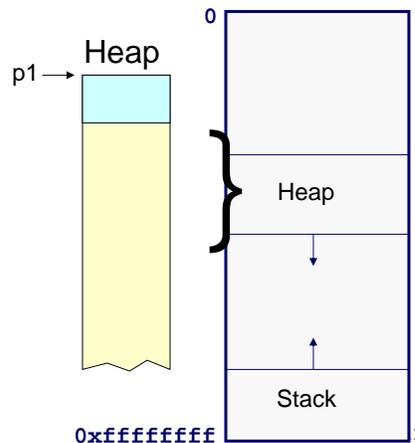
11

Heap: Dynamic Memory



```
#include <stdlib.h>
void *malloc(size_t size);
void free(void *ptr);
```

```
➔ char *p1 = malloc(3);
char *p2 = malloc(1);
char *p3 = malloc(4);
free(p2);
char *p4 = malloc(6);
free(p3);
char *p5 = malloc(2);
free(p1);
free(p4);
free(p5);
```

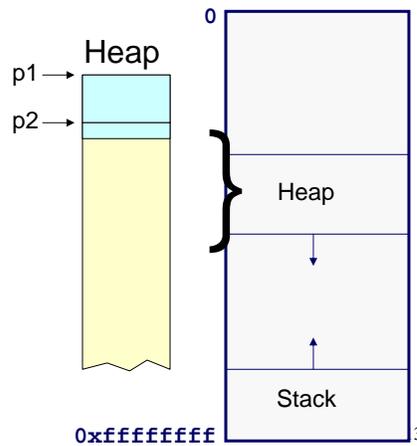


Heap: Dynamic Memory



```
#include <stdlib.h>
void *malloc(size_t size);
void free(void *ptr);
```

```
char *p1 = malloc(3);
➔ char *p2 = malloc(1);
char *p3 = malloc(4);
free(p2);
char *p4 = malloc(6);
free(p3);
char *p5 = malloc(2);
free(p1);
free(p4);
free(p5);
```

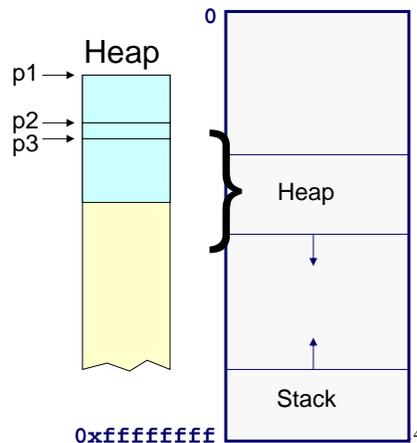


Heap: Dynamic Memory



```
#include <stdlib.h>
void *malloc(size_t size);
void free(void *ptr);
```

```
char *p1 = malloc(3);
char *p2 = malloc(1);
➔ char *p3 = malloc(4);
free(p2);
char *p4 = malloc(6);
free(p3);
char *p5 = malloc(2);
free(p1);
free(p4);
free(p5);
```

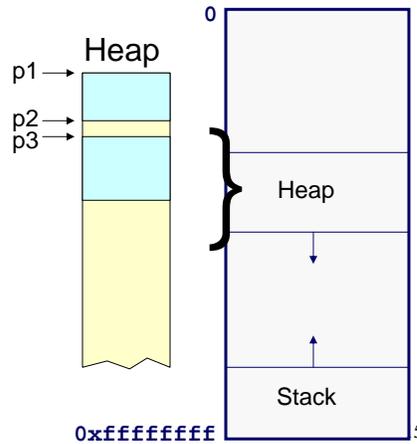


Heap: Dynamic Memory



```
#include <stdlib.h>
void *malloc(size_t size);
void free(void *ptr);
```

```
char *p1 = malloc(3);
char *p2 = malloc(1);
char *p3 = malloc(4);
➔ free(p2);
char *p4 = malloc(6);
free(p3);
char *p5 = malloc(2);
free(p1);
free(p4);
free(p5);
```

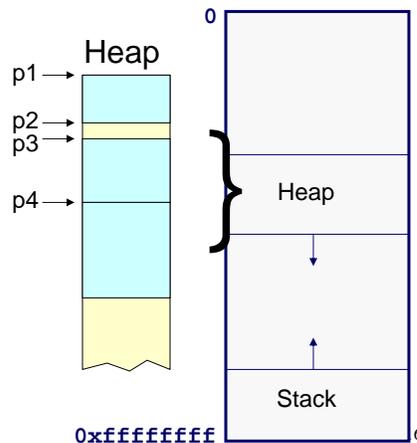


Heap: Dynamic Memory



```
#include <stdlib.h>
void *malloc(size_t size);
void free(void *ptr);
```

```
char *p1 = malloc(3);
char *p2 = malloc(1);
char *p3 = malloc(4);
free(p2);
➔ char *p4 = malloc(6);
free(p3);
char *p5 = malloc(2);
free(p1);
free(p4);
free(p5);
```

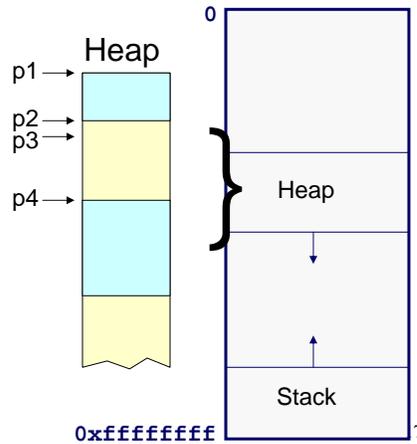


Heap: Dynamic Memory



```
#include <stdlib.h>
void *malloc(size_t size);
void free(void *ptr);
```

```
char *p1 = malloc(3);
char *p2 = malloc(1);
char *p3 = malloc(4);
free(p2);
char *p4 = malloc(6);
➔ free(p3);
char *p5 = malloc(2);
free(p1);
free(p4);
free(p5);
```

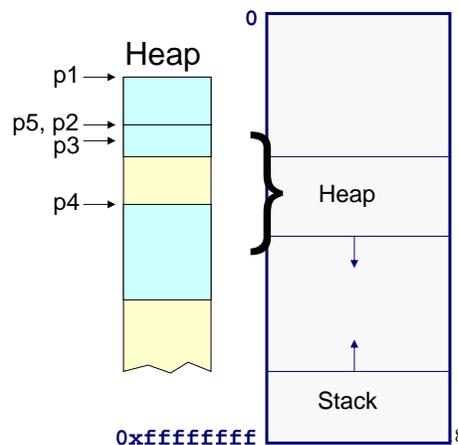


Heap: Dynamic Memory



```
#include <stdlib.h>
void *malloc(size_t size);
void free(void *ptr);
```

```
char *p1 = malloc(3);
char *p2 = malloc(1);
char *p3 = malloc(4);
free(p2);
char *p4 = malloc(6);
free(p3);
➔ char *p5 = malloc(2);
free(p1);
free(p4);
free(p5);
```

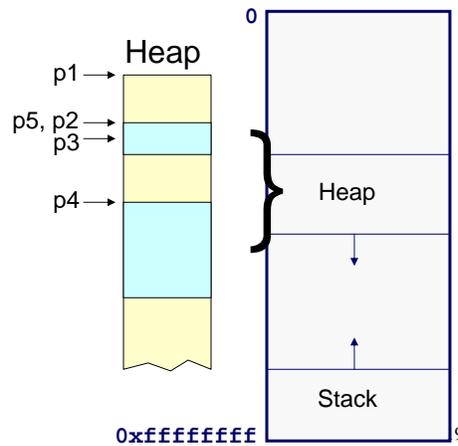


Heap: Dynamic Memory



```
#include <stdlib.h>
void *malloc(size_t size);
void free(void *ptr);
```

```
char *p1 = malloc(3);
char *p2 = malloc(1);
char *p3 = malloc(4);
free(p2);
char *p4 = malloc(6);
free(p3);
char *p5 = malloc(2);
➔ free(p1);
free(p4);
free(p5);
```

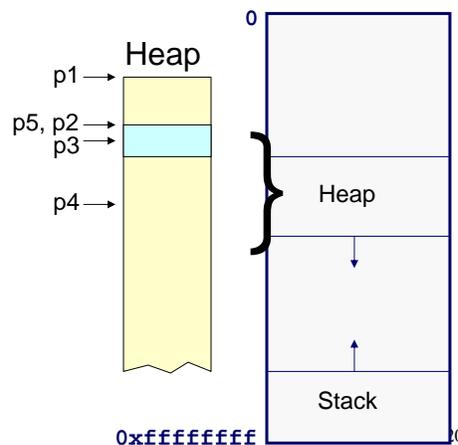


Heap: Dynamic Memory



```
#include <stdlib.h>
void *malloc(size_t size);
void free(void *ptr);
```

```
char *p1 = malloc(3);
char *p2 = malloc(1);
char *p3 = malloc(4);
free(p2);
char *p4 = malloc(6);
free(p3);
char *p5 = malloc(2);
➔ free(p4);
free(p5);
```

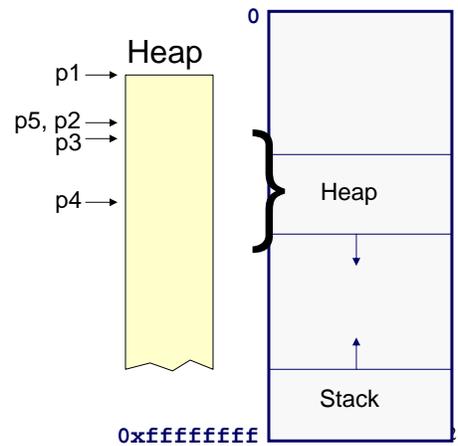


Heap: Dynamic Memory



```
#include <stdlib.h>
void *malloc(size_t size);
void free(void *ptr);
```

```
char *p1 = malloc(3);
char *p2 = malloc(1);
char *p3 = malloc(4);
free(p2);
char *p4 = malloc(6);
free(p3);
char *p5 = malloc(2);
free(p1);
free(p4);
➔ free(p5);
```



Goals for Malloc and Free



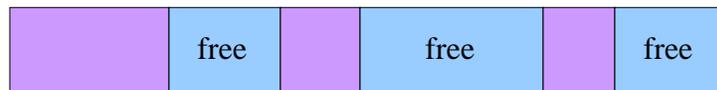
- Maximizing throughput
 - Maximize number of requests completed per unit time
 - Need both malloc() and free() to be fast
- Maximizing memory utilization
 - Minimize the amount of wasted memory
 - Need to minimize size of data structures
- Strawman #1: free() does nothing
 - Good throughput, but poor memory utilization
- Strawman #2: malloc() finds the “best fit”
 - Good memory utilization, but poor throughput

22

Keeping Track of Free Blocks



- **Maintain a collection of free blocks of memory**
 - Allocate memory from one of the blocks in the free list
 - Deallocate memory by returning the block to the free list
 - Ask the OS for additional block when more are needed
- **Design questions**
 - How to keep track of the free blocks in memory?
 - How to choose an appropriate free block to allocate?
 - What to do with the left-over space in a free block?
 - What to do with a block that has just been freed?

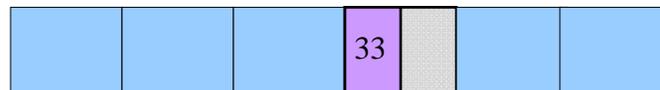


23

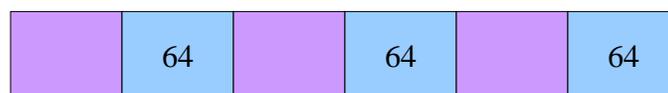
Need to Minimize Fragmentation



- **Internal fragmentation**
 - Allocated block is larger than malloc() requested
 - E.g., malloc() imposes a minimum size (e.g., 64 bytes)



- **External fragmentation**
 - Enough free memory exists, but no block is big enough
 - E.g., malloc() asks for 128 contiguous bytes



24

Simple “K&R-Like” Approach



- Memory allocated in multiples of a base size
 - E.g., 16 bytes, 32 bytes, 48 bytes, ...
- Linked list of free blocks
 - Malloc() and free() walk through the list to allocate and deallocate
- Malloc() allocates the *first* big-enough block
 - To avoid sequencing further through the list
- Malloc() *splits* the free block
 - To allocate what is needed, and leave the rest available
- Linked list is *circular*
 - To be able to continue where you left off
- Linked list in the order the blocks appear in memory
 - To be able to “coalesce” neighboring free blocks

25

Allocate Memory in Multiples of Base Size



- Allocate memory in multiples of a base size
 - To avoid maintaining very tiny free blocks
 - To align memory on size of largest data type (e.g., long)
- Requested size is “rounded up”
 - Allocation in units of `base_size`
 - Round: $(nbytes + base_size - 1) / base_size$
- Example:
 - Suppose `nbytes` is 37
 - And `base_size` is 16 bytes
 - Then $(37 + 16 - 1) / 16$ is $52 / 16$ which rounds down to 3



26

Linked List of Free Blocks



- Linked list of free blocks



- Malloc() allocates a big-enough block



- Free() adds newly-freed block to the list

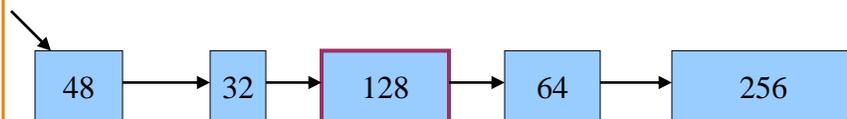


27

“First-Fit” Allocation



- Handling a request for memory (e.g., malloc)
 - Find a free block that satisfies the request
 - Must have a “size” that is big enough, or bigger
- Simplest approach: first fit
 - Sequence through the linked list
 - Stop upon encountering a “big enough” free block
- Example: request for 64 bytes
 - First-fit algorithm stops at the 128-byte block



28

Splitting an Oversized Free Block



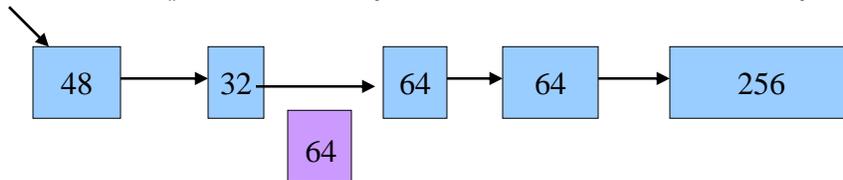
- Simple case: perfect fit

- Malloc() asks for 128 bytes, and free block has 128 bytes
- Simply remove the free block from the list



- Complex case: splitting the block

- Malloc() asks for 64 bytes, and free block has 128 bytes



29

Circular Linked List of Free Blocks

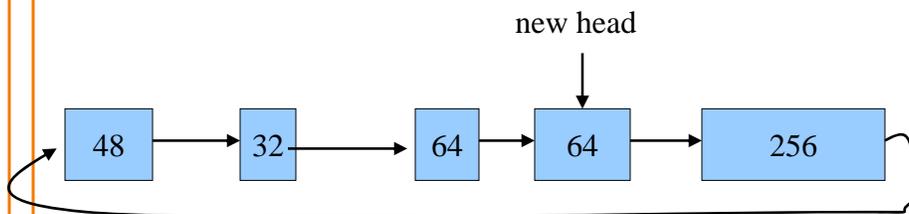


- Advantages of making free list a circular list

- Any element in the list can be the beginning
- Don't have to handle the "end" of the list as special

- Performance optimization

- Make the head be where last block was found
- More likely to find "big enough" blocks later in the list

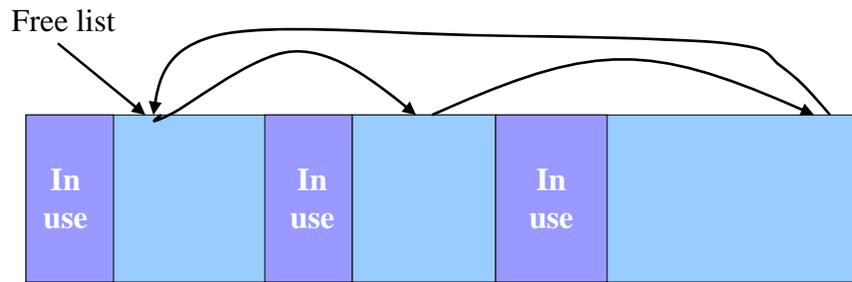


30

Maintaining Free Blocks in Order



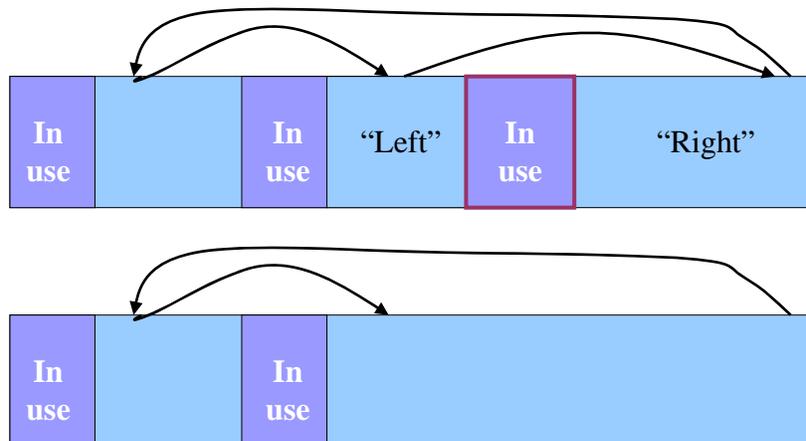
- Keep list in order of increasing addresses
 - Makes it easier to coalesce adjacent free blocks
- Though, makes calls to free() more expensive
 - Need to insert the newly-freed block in the right place



Coalescing Adjacent Free Blocks



- When inserting a block in the free list
 - “Look left” and “look right” for neighboring free blocks



An Implementation Challenge



- Need information about each free block
 - Starting address of the block of memory
 - Length of the free block
 - Pointer to the next block in the free list
- Where should this information be stored?
 - Number of free blocks is not known in advance
 - So, need to store the information on the *heap*
- But, wait, this code is what *manages* the heap!!!
 - Can't call `malloc()` to allocate storage for this information
 - Can't call `free()` to relinquish the storage, either

33

Store Information in the Free Block



- Store the information directly in the free block
 - Since the memory isn't being used for anything anyway
 - And allows data structure to grow and shrink as needed
- Every free block has a header
 - Size of the free block
 - Pointer to (i.e., address of) the next free block



- Challenge: programming outside the type system

34

Conclusions



- **Elegant simplicity of K&R malloc and free**
 - Simple header with pointer and size in each free block
 - Simple circular linked list of free blocks
 - Relatively small amount of code (~25 lines each)
- **Limitations of K&R functions in terms of efficiency**
 - Malloc requires scanning the free list
 - To find the first free block that is big enough
 - Free requires scanning the free list
 - To find the location to insert the to-be-freed block